

构建以事件为核心的长期保存系统起源管理框架

吴振新¹ 李文燕² 蒋世银¹

¹中国科学院文献情报中心 北京 100190

²中新金桥数字科技(北京)有限公司 北京 100096

摘要: [目的/意义]研究建立长期保存系统起源管理框架,通过有效管理起源信息,确保长期保存系统所存档数据的真实可靠可用。[方法/过程]基于数字对象保存周期进行起源事件定义,基于 OAIS 保存流程进行起源管理框架设计,以事件为核心进行起源管理功能模型和起源信息模型设计。[结果/结论]初步完成基于事件的保存系统起源管理框架的设计,既遵循了保存领域的相关标准,同时兼顾了实践需求,对长期保存系统具有很好的普适性和可行性,但其在有效性和实用性方面还有待进一步验证。

关键词: 起源信息 起源事件 长期保存 生命周期管理 OAIS PREMIS

PROV 语义模型

分类号: G250.76

前言

数字对象的起源信息,即 provenance,它记录了数字对象的变化历史。通过起源信息,人们可以全面了解数字对象产生之后所发生的变化以及变化的原因、时间、地点、相关人员等 7W 信息(what、where、who、when、which、why、how)。

数字资源长期保存系统作为一类特殊的数据管理系统,通过摄入、保存、管理等一系列管理行为,确保数字对象经过足够长时间后还能够被指定社团所使用,对于数据的真实可靠可用,其面临着更大的责任和挑战。在长期保存系统中,起源信息能够发挥多方面的作用。一方面长期保存系统要在相当长的时间内管理和保持数字对象的可用性,另一方面它要抵抗技术变化对数字对象及保存系统所带来的影响,格式迁移、媒体迁移、技术更新是长期保存的常用策略,所以无论当对象本身发生变化还是其所处环境发生改变,都可以利用起源详细记录这些改变,并维护这些变化前后的数字对象的关联。通过这种方法,使得保存系统能够有效进行版本和衍生物管理,并为数字对象的真实性和系统的可信赖认证提供证据,同时还能权限管理和责任归属提供支持。因此起源对于保存系统有着更为重要的意义。

本文作者曾在《起源技术在长期保存中的研究与应用》^[1]一文中全面总结和分析了起源在长期保存中的研究情况,初步提出了一个起源管理框架,本文将基于该文的基础,概述如何进一步完善长期保存系统起源管理框架和相关功能的设计。

作者简介: 吴振新(ORCID: 0000-0003-4966-1961), 研究馆员, wuzx@mail.las.ac.cn; 李文燕(ORCID: 0000-0002-7695-5087), 助理工程师, 硕士。蒋世银(orcid: 0000-0002-2038-4027), 助理馆员, 硕士。

1 长期保存系统的起源管理框架基本设计思路

OAIS 作为长期保存领域的基础标准, 已经成为保存系统的基准参考。在 OAIS 中, 起源被定义为内容信息的历史, 展示了内容信息从产生以后发生的相关变动^[2], 并将其作为保存描述信息 (Preservation Description Information, p;PDI) 的一部分。起源信息的管理贯穿了数字对象在 OAIS 系统中的整个生命周期。考虑诸多因素, 本文在起源管理框架的设计中, 遵循了如下的设计原则:

(1) 基于 OAIS。OAIS 是长期保存的通用标准, 提供了长期保存的基本流程和事件, 所以它是本文研究框架的最基本的出发点。

(2) 基于数字对象的保存周期。本文以数字对象提交到保存系统为起点, 对其进入保存系统后的整个保存周期内所有变化来实施起源采集与管理。在摄入之前的起源信息, 可由内容生产者在和保存方协商一致的基础上以规范的方式提交到长期保存系统, 这一点和 OAIS 的要求也是一致的。

(3) 以事件为核心来记录起源信息。在长期保存的过程中, 数字对象会因各种管理活动产生多种事件并产生起源信息, 可以说事件通常都伴随着起源信息的产生。

(3) 交互性。信息模型是系统软件设计中重要的内容, 为增强起源信息在不同系统之间的交互性, 起源管理应采用信息模型来组织管理数据。

(4) 通用性。本框架旨在为长期保存中组织和管理起源信息提供一般的功能流程、模型等内容的参考, 与具体技术实现无关。

2 面向保存周期管理的起源事件定义

事件是对象发生变化的主要驱动, 通过事件可以把各种类型的状态变化串联起来。随着事件的累积, 发生在对象上的事件链可以动态的呈现保存对象的状态改变。因此本文的起源事件定义为涉及或影响至少一个对象或代理的保存系统可识别的动作, 和计算机中常说的单击、双击、窗体加载等事件是不同的, 它是长期保存系统定义的保存系统处理对象的操作或操作集, 如压缩文件、摄入信息包、创建对象等。识别记录哪些事件作为起源事件是本文起源管理框架的核心问题。

PREMIS^[3]作为保存领域的保存元数据标准, 其定义了五个基本实体, 事件是其中之一, 因此使用事件记录起源可以便于使用长期保存元数据进行描述。PREMIS 中定义了 15 种保存事件: creation、deaccession、decompressionn、decryption、deletion、digital signature validation、dissemination、fixity check、ingestion、message digest calculation、migration、normalization、replication、validation 和 virus check。但这些事件不是专门针对起源设计的起源事件, 同时不同事件之间还有重叠的地方, 如 creation 和 normalization; 部分事件指向不明, 在事件应用中可能会产生歧义, 如 validation。所以在实际应用中, 需要明确定义长期保存系统中哪些事件应该被记录为起源事件。

OAIS 认为: 起源是内容信息的历史, 它展示了内容信息产生的由来, 从产生以后发生的变化, 以及自创建以后的保管责任方的改变。这个定义了暗含选择起源事件遴选两个重要依据, 即“时间”和“变化”。“时间”, 即数字对象的保存周期, 在 OAIS 中整个保存周期包含了 6 个保存流程, 即摄入 (Ingest)、归档存储 (Archival Storage)、数据管理 (Data Management)、业务管理

(Administration)、保存规划 (Preservation Planning) 和访问 (Access)。
“变化”，即判断一个事件是否为起源事件的依据。

综上所述，在遴选起源事件时应考虑如下方面：

- 导致内容对象最初产生，这是一个由无到有的过程。
- 导致内容对象本身发生变化，或者能够捕获长期保存过程中，自然发生的变化，这些变化包括内容、结构、数量、格式、位置、元数据和保管责任方。
- 导致新版本对象的产生，虽然数字对象内容本身并未发生任何变化，这是长期保存基本要义，但是产生了和该数字对象关系密切的新对象，例如副本，不同格式的新版本，对于对象的复用都很有益处。
- 导致数字对象的版权和管理权发生变化。
- 导致数字对象的消亡。

根据以上原则，从 OAIS 所包括的流程中遴选出如表 1 所示的起源事件。

表 1 长期保存起源事件清单

事件	英文名称	涉及流程
捕获	capture	摄入
解密	decryption	摄入
逆压缩	decompression	摄入
压缩	compression	摄入
病毒检查	virus check	摄入、归档存储
数量检查	number check	摄入、归档存储
内容检查	content check	摄入、归档存储
不变性检查	fixity check	摄入、归档存储
备份检查	backup check	归档存储
硬件检测	hardware detection	摄入、归档存储
格式检查	format check	摄入、归档存储
目录删除	deaccession	归档存储
生产者登记	producer register	摄入
摄入消息计算	message digest calculation	摄入
规范化	normalization	摄入
当前化	contemporary	摄入、归档存储
元数据抽取	description	摄入
创建	creation	摄入
摄入	ingestion	摄入
备份	replication	摄入、归档存储
媒体迁移	media migration	归档存储
删除	delete	归档存储
数据恢复	dataRecovery	归档存储
模式更新	schema update	数据管理
分发	dissemination	访问

数据迁移	transfer	访问
------	----------	----

3 嵌入 OAIS 保存流程的起源管理框架

从起源事件清单可以看出，起源事件涉及 OAIS 的所有流程和功能模块，起源信息的管理需要嵌入到保存系统的完整流程中，如图 1 所示。

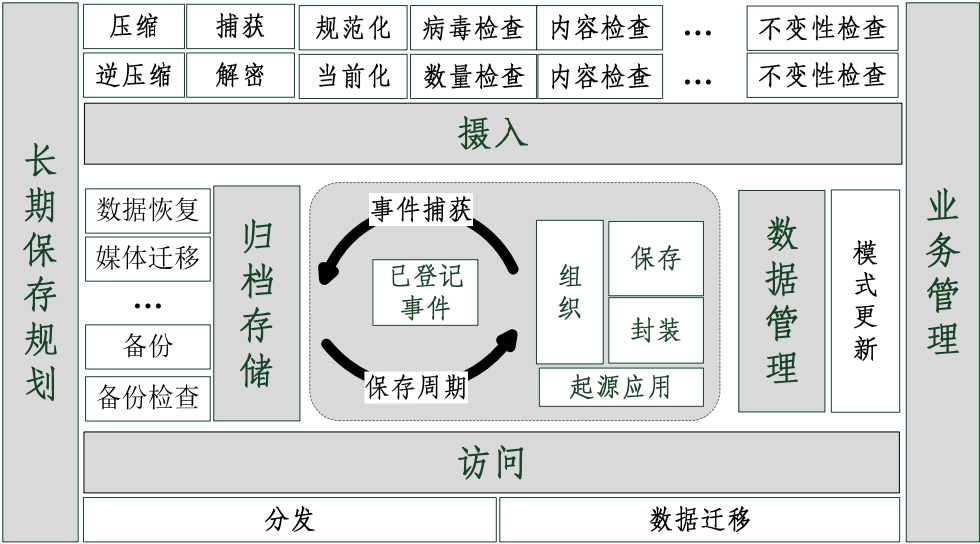


图 1. 嵌入 OAIS 保存流程中的起源管理框架

其中，图 1 中心部分是起源的管理模块，它要嵌入 OAIS 的各个流程中，动态地监测各个流程的事件，并根据起源管理中预先配置好的起源事件清单，捕获数字对象的起源事件。然后把长期保存捕获的事件和生产者提供的事件按照相应的起源模型组织成规范的起源，存储为相应的格式（封装与存储），并由保存管理模块对组织模块生成的起源进行长期保存和管理，保证起源的完整性、可理解性和长期可访问性，同时应用模块按保存系统要求提供为用户或长期保存的其他模块（如真实性管理）提供起源的使用。

4 以事件为核心的起源管理功能模型

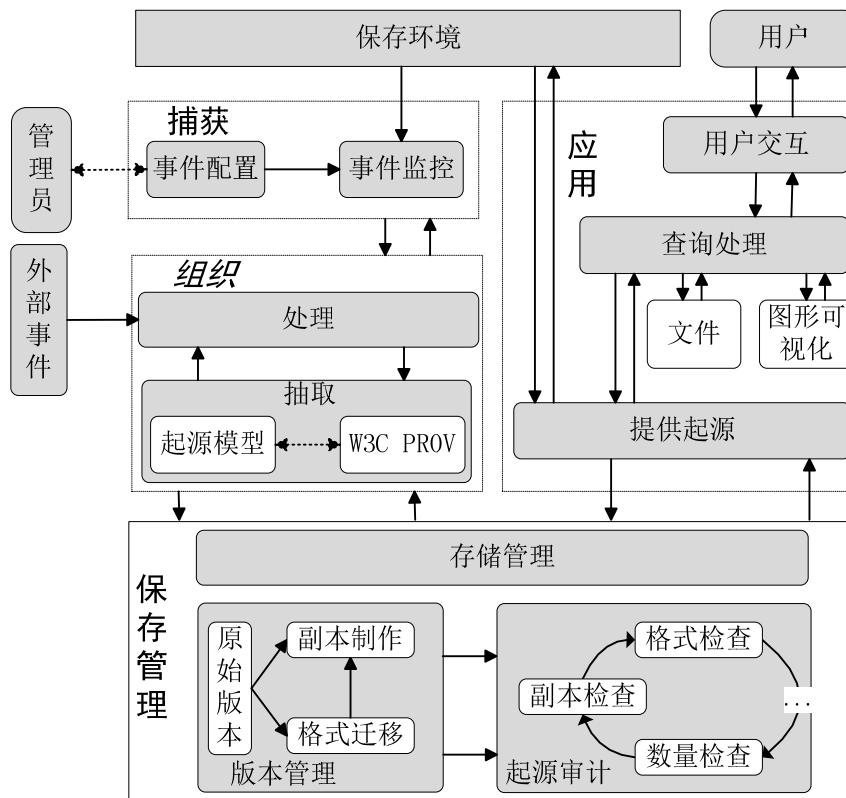


图 2. 以事件为核心的起源管理功能模型

图 2 清晰的展示了起源管理各个功能模块的相互关系及数据流向。该模型包含了四个基本子功能模块即捕获、组织、保存管理和应用，其中捕获、组织和保存管理是功能模型的重点。

（1）捕获模块

事件配置功能负责预定义和配置起源管理需要捕获的事件类型。这个功能在捕获之前完成，由长期保存系统的管理人员根据保存系统的功能所包含的各个操作，归纳出需要记录为起源的事件，对其进行详细的定义，并把起源事件清单配置为计算机可读的格式，如数据库表或者 XML 文件。

事件监控功能负责动态地监测保存系统所发生的所有事件，当某个事件和预定义起源事件清单中事件相匹配时，则触发组织模块，把事件信息如事件内容、事件时间、操作对象和使用设备等内容传递给组织模块。

（2）组织模块

组织模块负责根据接收到的事件消息后将其添加到起源记录的任务队列中，供抽取功能使用。通过这种异步记录的方式来组织起源，既能减少对系统原有进度的影响，又能减少服务器的负担。起源记录的任务队列包括两种类型事件，一种是自动捕获发生在保存系统内部的事件，另一种是由被保存内容的生产者的提供的外部事件。

抽取功能按照顺序读取任务队列的事件信息，根据系统设定的信息模型（如 xml schema）对事件信息进行规范化组织，生成规范化起源。

（3）保存管理模块

存储管理功能把从组织模块接收到的起源按照相关的方式进行存储，并维护起源数字对象之间的关联关系。版本管理功能实现对起源的各个版本的管理。主要是按照保存计划 and 政策进行副本制作或者支持格式迁移。起源审计功能负

责为每个版本的起源定期执行不变性检查（Fixity Check）、格式检查和副本检查。触发该功能有两种类型任务，一种是定期检查任务，一种是新增加起源、备份起源和改变起源版本触发的任务。

(4) 应用模块

应用模块为系统的其他模块（如审计追踪）使用起源信息提供标准的接口调用。查询处理功能直接接收用户的起源请求，然后调用相关的起源接口，返回特定格式的起源给用户。用户交互功能则为用户提供可视化界面，如网页。起源查询或下载的请求消息由用户在交互界面中发起后被传递给查询处理模块，将底层处理后从交互界面返回用户请求的特定格式起源。

5 以事件为核心的起源信息模型

使用信息模型不仅对所管理的数据进行有效组织，同时也有利于长期保存、管理和重用。OAIS 中数字对象包括内容信息和表征信息两部分，因此起源不但要记录内容信息的变化，也同时需要记录表征信息的变化。

因此起源信息中该包含以下相关内容：

(1) 事件

事件的驱动使得数字对象发生变化。事件的细节描述是起源信息的重点内容，除了包含事件标识符、细节描述、时间、事件类型、处理设备、处理结果、地点和发生原因，还要包含事件涉及的责任人和被操作的对象基本信息。

(2) 数字对象

需要完整记录事件涉及到的数字对象。通过在事件中引用数字对象的标识符将二者关联起来，但不包含数字对象的描述元数据。如果一个事件同时关联两个甚至更多数字对象，这就意味着所有对象都拥有该起源信息，应该包含所有的数字对象的标识符。

(3) 代理内容

狭义的代理（Agent）指的是事件的操作人，此处的代理的含义更广，它包括组织、个人、软件和硬件 4 种内容。

(4) 数字对象之间的关系

事件对数字对象的操作可能会导致新版本对象产生，如副本或不同格式的副本。虽然数字对象之间的关系可能不被直接记录在起源信息中，但是通过分析相关事件的性质和涉及的输入、输出对象可以间接得出数字对象的版本变化。

W7 语义模型^[4]给出了从 7 个维度来记录起源信息的思路，较全面地说明了起源信息包含的内容，对于构建起源模型具有重要的参考作用。本文在 W7 语义模型和上文归纳的事件起源内容基础上设计如图 3 所示的起源信息模型，此模型从事件的角度描述了起源信息应该包含的内容元素，并从 7 个维度来设计每个起源事件包含的内容概念。

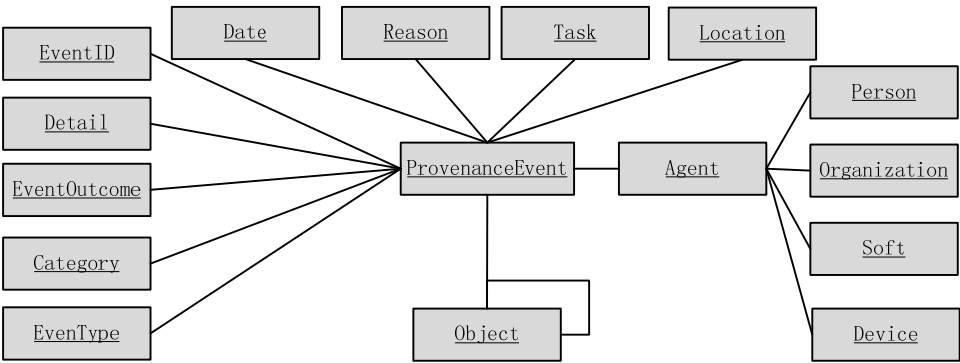


图 3 事件起源信息模型

在该信息模型中,以事件为中心把记录数字对象变化的各种信息串联起来,涵盖了起源事件的基本元素: Object、Agent、EventID、Date、Reason、Task、Detail、EventOutcome、Category、 EventType 和 Location, 每一个元素都能对应表示 W7 模型的一个维度, 如表 3 所示。虽然 Agent 和 Object 是事件的一部分, 但对它们详细的描述元数据设计不包含在框架模型的范围内, 该模型只对起源事件应该包含的各个概念加以描述。

表 2 事件起源信息模型基本元素表

基本元素	解释	对应维度
EventID	事件的唯一标识符, 用以引用该事件。	What
Date	记录事件发生的时间, 可能是一个时间点, 也可能是一个时间段。	When
Detail	事件的详细内容的文字描述, 例如“成功执行对图片的不变性检查”。	What
EventOutcome	事件的执行的结果描述。	What
EventType	事件类型, 如病毒检查、不变形检查, 推荐使用受控词汇。	What
Category	事件分类, 如病毒检查和不变形验证均属于验证类事件, 对不同事件进行归类, 方便查询。	What
Task	驱动事件发生的任务技术, 类似于工作流的定义。	How
Reason	驱动事件发生的原因描述。	Why
Location	和事件相关的位置信息。	Where
Agent	和事件相关的代理, 包括 4 种子类型, Person (人)、Organization (组织)、Soft (软件) 和 Device (物理设备)。	Who, Which
Object	事件涉及的被保存对象。	What

作者在系统实现时复用了 PREMIS OWL^[5]和 W3C PROV-O^[6]来实现起源组织模型, 并利用 RDF 进行了起源封装, 同时本文没有涉及起源的存储和封装策略, 这些内容将在作者的另一篇论文中详细介绍。

6 结语

目前该框架系统原型刚刚完成, 其有效和实用性还有待进一步验证。总的来看, 基于事件的保存系统起源管理框架的设计, 不但遵循了保存领域的相关标准, 同时也是从大量的项目实践的调研中提炼、总结基础上完成的^[1], 兼顾了理论和实践两个方面, 对长期保存的信息系统具有很好的普适性和可行性。该框架基于数字对象的保存周期, 提出了嵌入到 OAIS 流程中与相关保存事件相融合的方法, 保证了起源的有效和完整。而以事件来记录起源的方式可以有效的把分散的元数据碎片按照时间组织起来, 既紧抓起源产生的本质, 又便于起源信息采集管理, 这一点和 PREMIS 认为起源以事件为驱动的观点是一致的。另外本框架的起源信息组织采用了信息模型的设计, 与具体的实现技术无关, 具有高度的抽象性, 可以适合大多数长期保存系统。

保存领域的研究人员已经充分认识到起源对长期保存的重要作用,不断探索对起源的有效管理和使用,希望本文所做的研究和努力,能够为相关人员在长期保存的起源管理理论和实践方面提供有益的参考。

-
- [1]吴振新,李文燕. 起源技术在长期保存中的应用与研究[J]. 图书情报工作, 2015, 59(8):118-125.
- [2]CCSDS 650.0-M-2, Reference Model For An Open Archival Information System (OAIS)[S].
- [3] PREMIS Editorial Committee. PREMIS data dictionary for preservation metadata, version 2.0[J]. Retrieved November, 2008, 15: 2009.
- [4]Ram S, Liu J. A Semantic Foundation for Provenance Management[J]. Journal on Data Semantics, 2012, 1(1):11-17.
- [5]PREMIS Editorial Committee. PREMIS OWL ontology 2.2 now available[EB/OL]. [2015-03-31]. <http://www.loc.gov/standards/premis/ontology-announcement.html>
- [6]PROV-O: The PROV Ontology [EB/OL]. [2015-5-17].<https://www.w3.org/TR/2013/REC-prov-o-20130430/>

作者贡献说明:吴振新:提出研究思路,设计论文框架,论文修改、定稿;李文燕:论文资料查询、收集,论文初稿撰写。蒋世银:论文资料查询、收集。

Construct a Provenance Management Framework for Long Term Preservation System Based on the Event-Centric Method

Wu Zhenxin¹ Li Wenyan^{1,2} Jiang Shiyin¹

¹National Science Library, Chinese Academy of Sciences, Beijing 100190

²KingChannels Co. Ltd, Beijing 100196

Abstract:

[Purpose/Significance] Research on constructing a provenance management framework for long-term preservation system, by providing the effective management of provenance information, we would ensure the authenticity, reliability and usability of data objects which are archived in the long-term preservation system. **[Method/Process]** Complete the definition of provenance events based on preservation lifecycle, complete the design of a provenance management framework according to OAIS, do the design of the management function model and information model by the way of event-centric. **[Result/Conclusion]** Construct a provenance management framework for long-term preservation system, not only follow the relevant standards of long-term preservation, but also take into account the needs of practice, although its effectiveness and practicability remains to be verified furthermore, it still has good applicability and feasibility for the long term preservation system.

Keywords: Provenance Information Provenance Event Long-term Preservation Lifecycle Management OAIS PREMIS PROV Ontology